# Time Series Analysis

**Neychev Radoslav**

ML Instructor, BigData Team

Research Scientist, MIPT
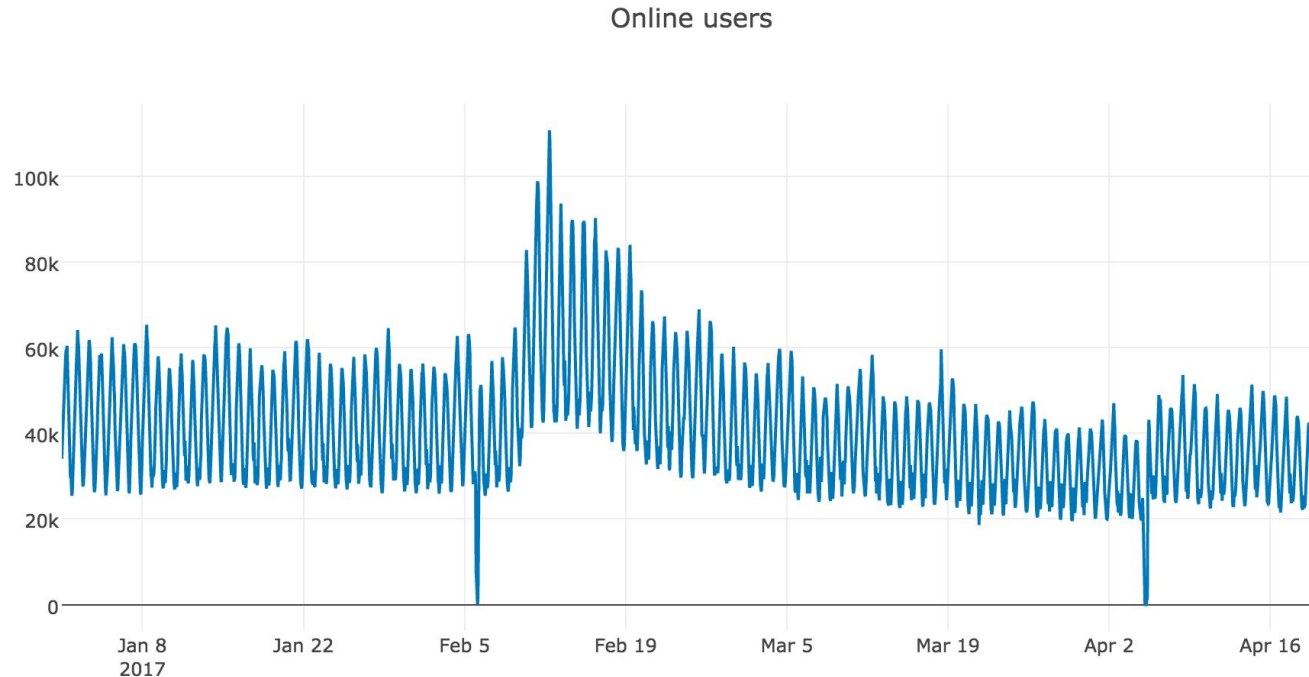
27.08.2018, Moscow, Russia

1.  Time Series as data with time hierarchy.

2.  Classical approaches: AR, MA, ARIMA, exp. smoothing etc.

3.  Practice.

4.  Feature engineering for Time Series analysis.

5.  Modern approaches.

6.  (extra) Neural networks (CNN & RNN) in TS analysis.

7.  More practice.

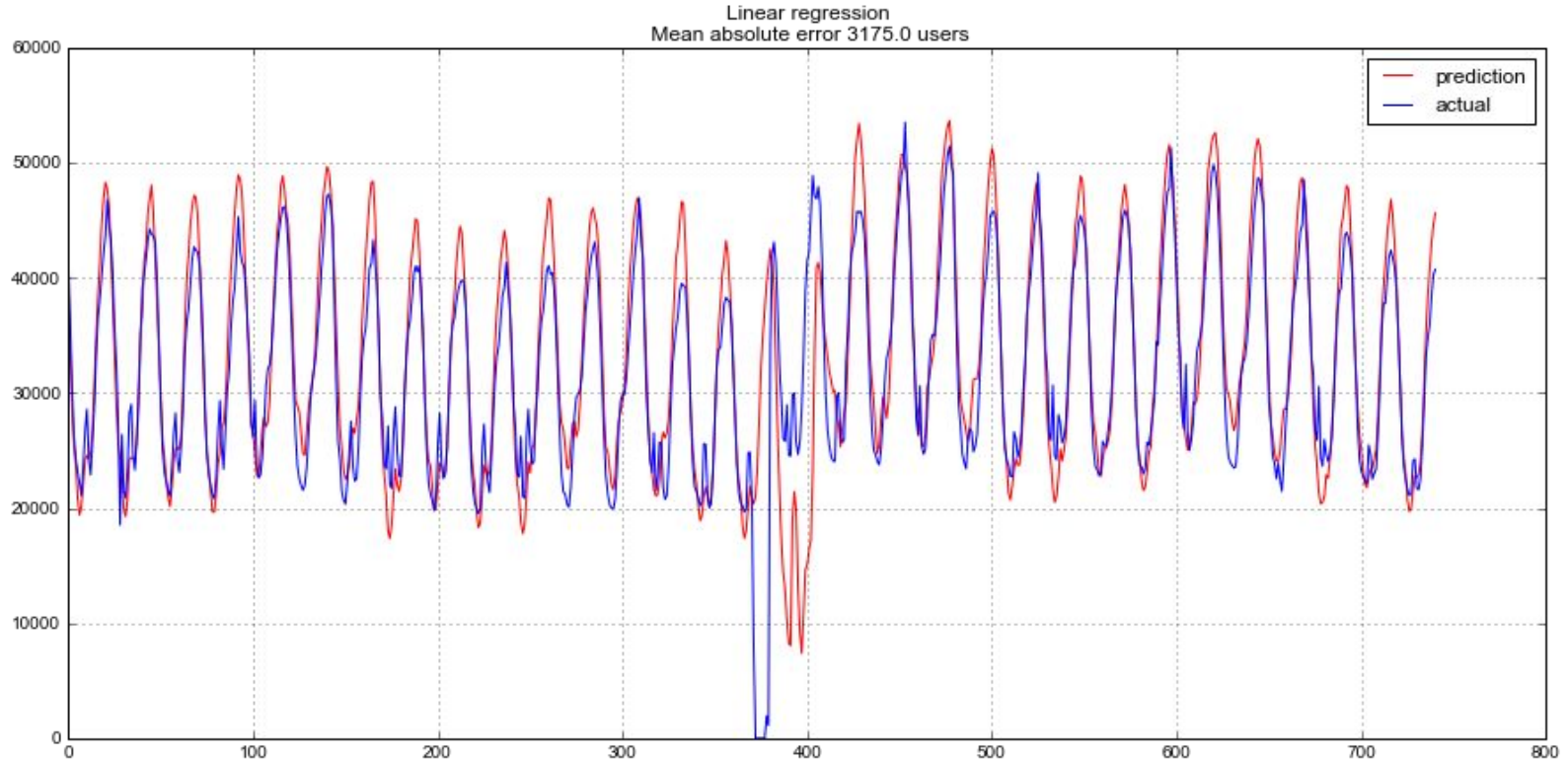Materials: http://bit.ly/ml4megafon_august18_public

Credits: these slides are based on ODS open ML course materials and several blog posts.

*Time Series* -- series of data observations with time hierarchy.

Online users

**BIGDATA TEAM**

How to measure the prediction quality?



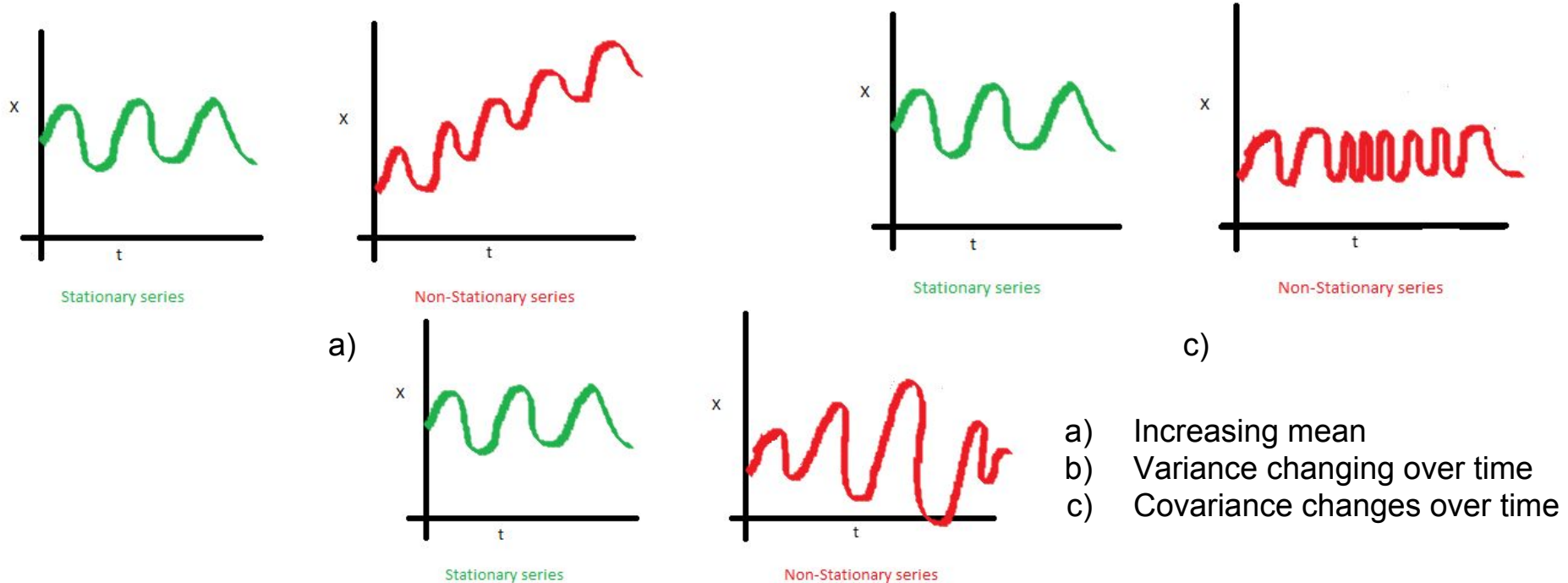Linear regression
Mean absolute error 3175.0 users

# Time Series: metrics in forecasting

- [R squared](#) (coefficient of determination): $R^2(y, \hat{y}) = 1 - \dfrac{\sum_{i=0}^{n_{\text{samples}}-1}(y_i - \hat{y}_i)^2}{\sum_{i=0}^{n_{\text{samples}}-1}(y_i - \bar{y})^2}$

- [Mean Absolute Error](#): $\text{MAE}(y, \hat{y}) = \dfrac{1}{n_{\text{samples}}} \sum_{i=0}^{n_{\text{samples}}-1} |y_i - \hat{y}_i|.$

- [Median Absolute Error](#): $\text{MedAE}(y, \hat{y}) = \text{median}(|\, y_1 - \hat{y}_1 \,|, \ldots, |\, y_n - \hat{y}_n \,|).$

- [Mean Squared Error](#): $\text{MSE}(y, \hat{y}) = \dfrac{1}{n_{\text{samples}}} \sum_{i=0}^{n_{\text{samples}}-1} (y_i - \hat{y}_i)^2.$

- Mean Absolute Percentage Error: $\text{MAPE} = \dfrac{1}{n_{\text{samples}}} \sum_{i=0}^{n_{\text{samples}}-1} \dfrac{|y_i - \hat{y}_i|}{|y_i|}$

- [Mean Squared Logarithmic Error](#): $\text{MSLE}(y, \hat{y}) = \dfrac{1}{n_{\text{samples}}} \sum_{i=0}^{n_{\text{samples}}-1} (\log_e(1 + y_i) - \log_e(1 + \hat{y}_i))^2.$

*Stationary process* does not change its statistical properties (mean and variance) over time.



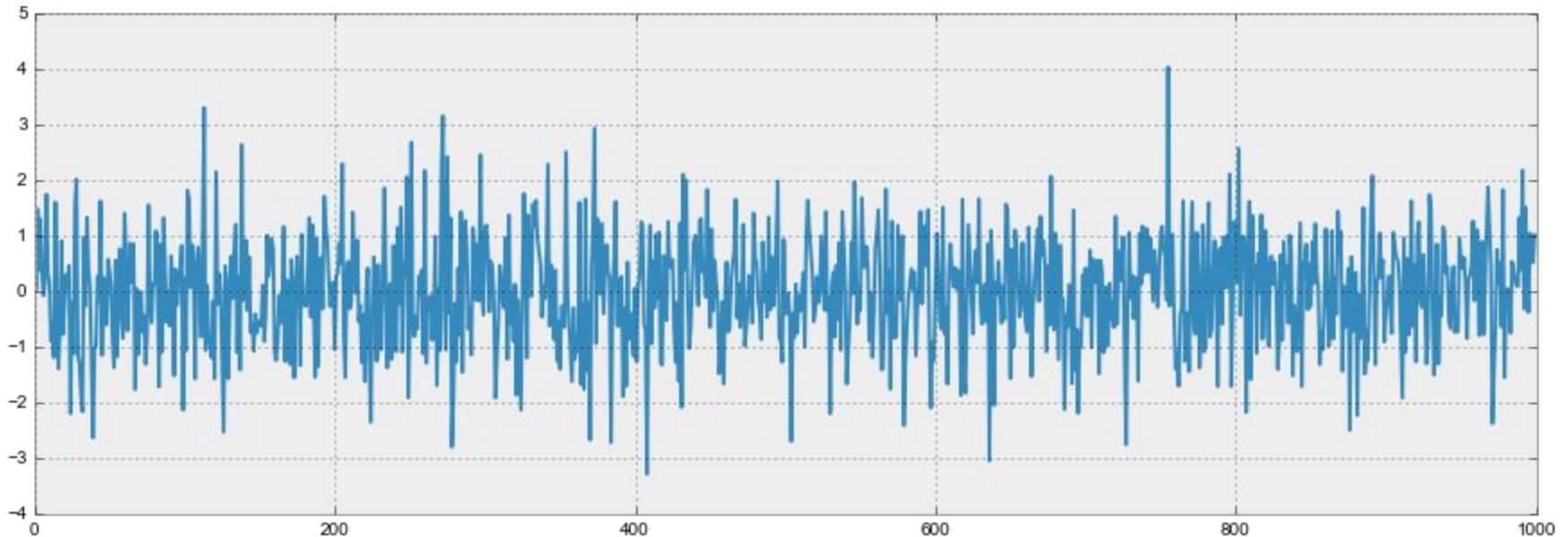Stationary series | Non-Stationary series

a)

Stationary series

Stationary series | Non-Stationary series

c)

b)

a) Increasing mean
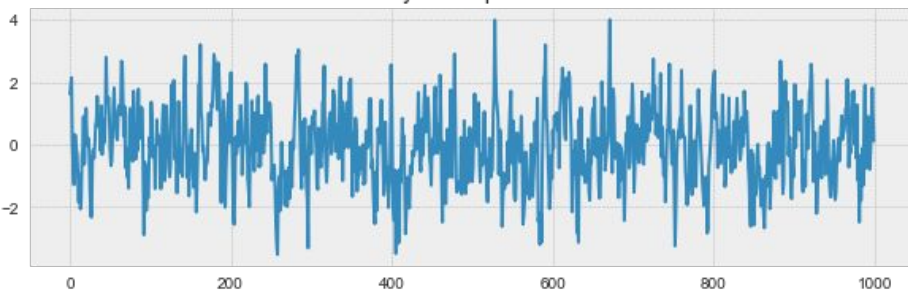b) Variance changing over time
c) Covariance changes over time

Source: post by Sean Abu
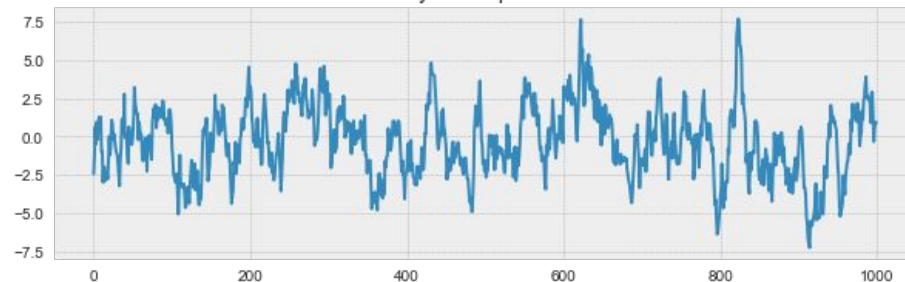
Process oscillates around 0 with deviation of 1.

Let's transform the white noise process: $y(t) = \rho * y(t-1) + e(t)$



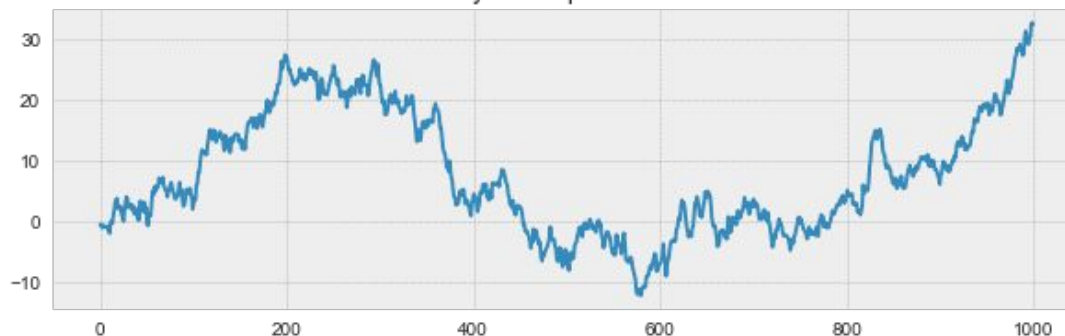Augmented Dickey-Fuller test in statsmodels: [link](#)
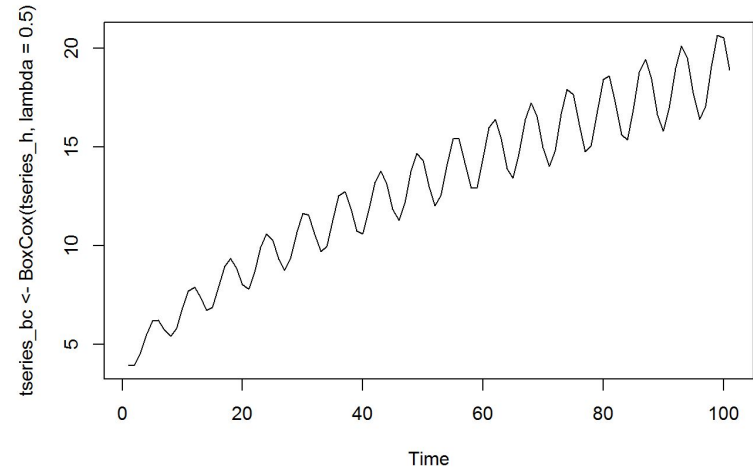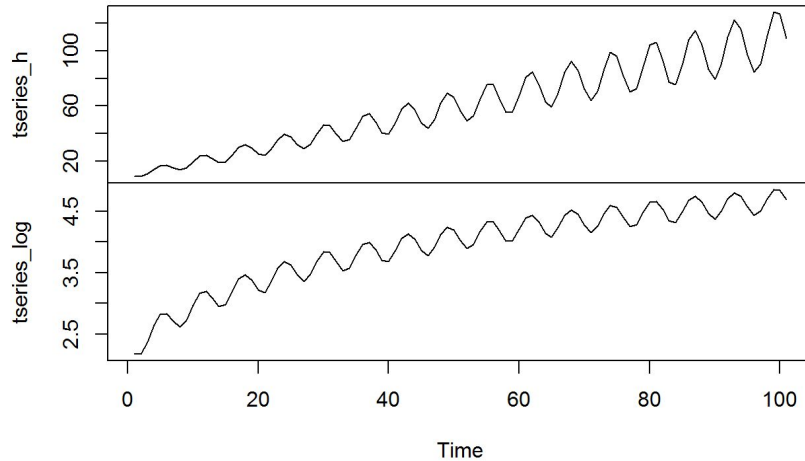
Box-Cox transformation:

$$w_t = \begin{cases} \log y_t, & \text{if } \lambda = 0; \\ \dfrac{(y_t^{\lambda}-1)}{\lambda}, & \text{otherwise} \end{cases}$$
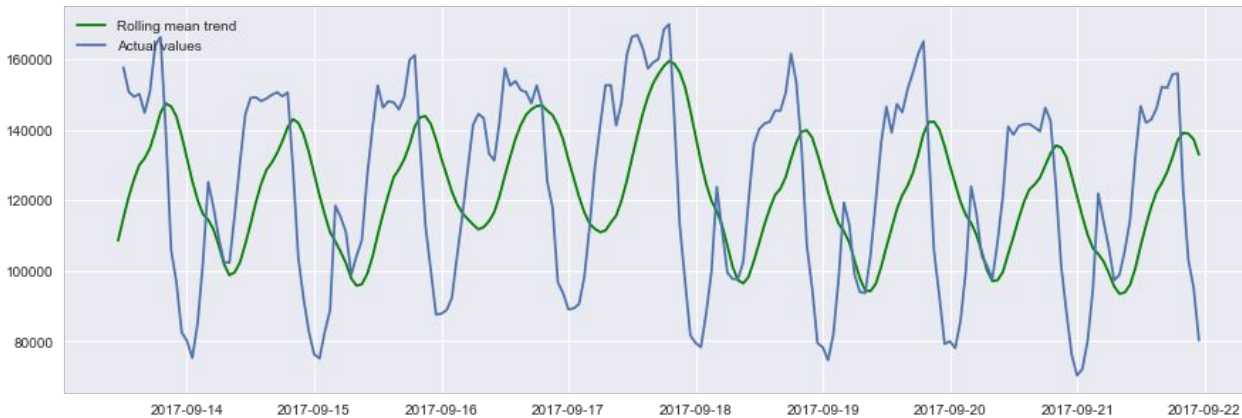


Source: datascienceplus blog post

- Moving average: $\hat{y}_t = \frac{1}{k} \sum_{n=0}^{k-1} y_{t-n}$

- Weighted moving average: $\hat{y}_t = \sum_{n=1}^{k} \omega_n y_{t+1-n}$

- Exponential smoothing: $\hat{y}_t = \alpha \cdot y_t + (1 - \alpha) \cdot \hat{y}_{t-1}$
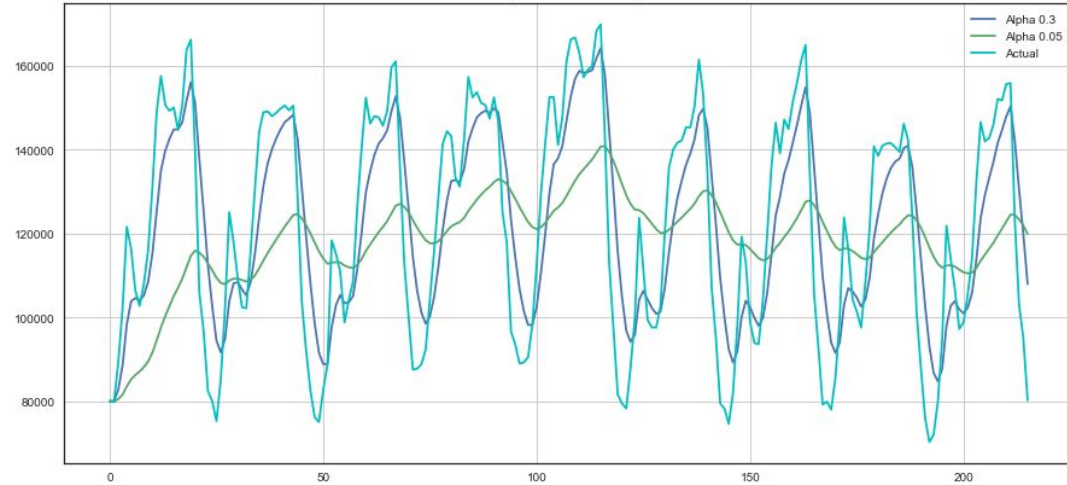
Caution: only one step is predicted

Exponential Smoothing

Moving average
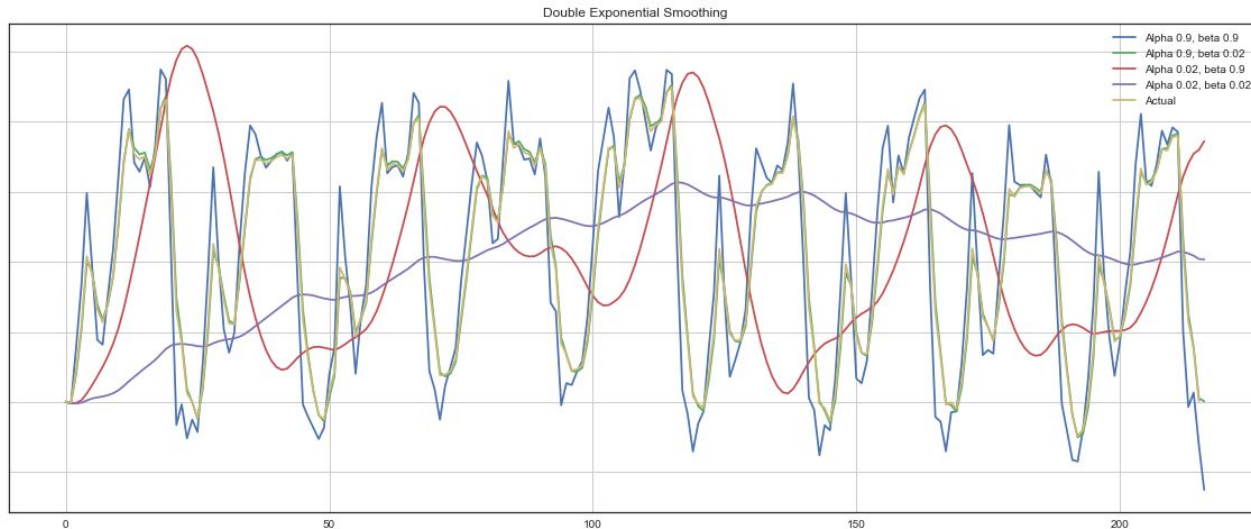window size = 12

# Double exponential smoothing

Two components:

intercept (or level) :

$$\ell_x = \alpha y_x + (1 - \alpha)(\ell_{x-1} + b_{x-1})$$

trend (or slope):

$$b_x = \beta(\ell_x - \ell_{x-1}) + (1 - \beta)b_{x-1}$$

$$\hat{y}_{x+1} = \ell_x + b_x$$



Double Exponential Smoothing

Now two steps can be predicted

Let's add one more component: seasonality (s).

Seasonal component in the model should explain repeated variations around intercept and trend, and it will be described by the length of the season.

$$\ell_x = \alpha(y_x - s_{x-L}) + (1 - \alpha)(\ell_{x-1} + b_{x-1})$$

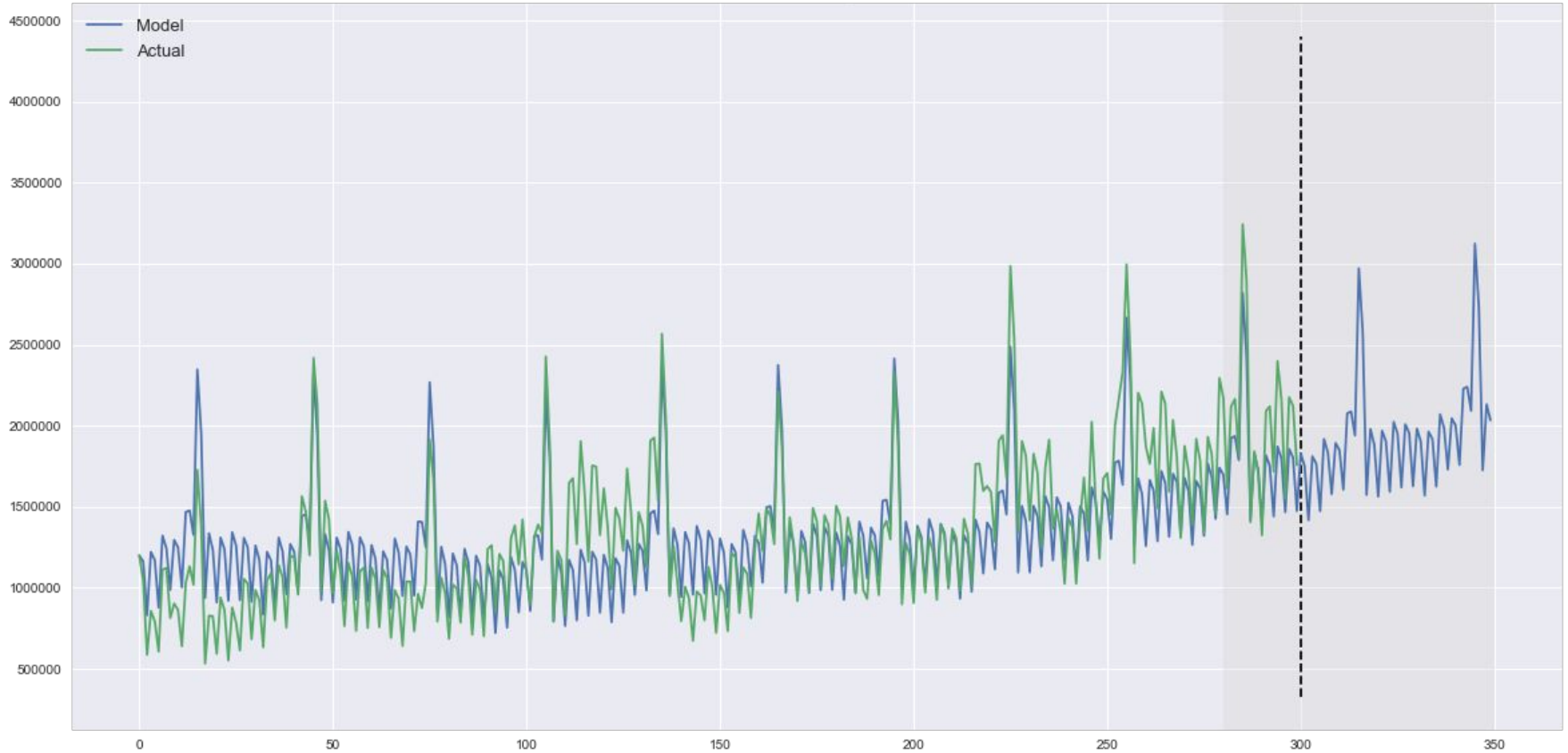$$b_x = \beta(\ell_x - \ell_{x-1}) + (1 - \beta)b_{x-1}$$

$$s_x = \gamma(y_x - \ell_x) + (1 - \gamma)s_{x-L}$$
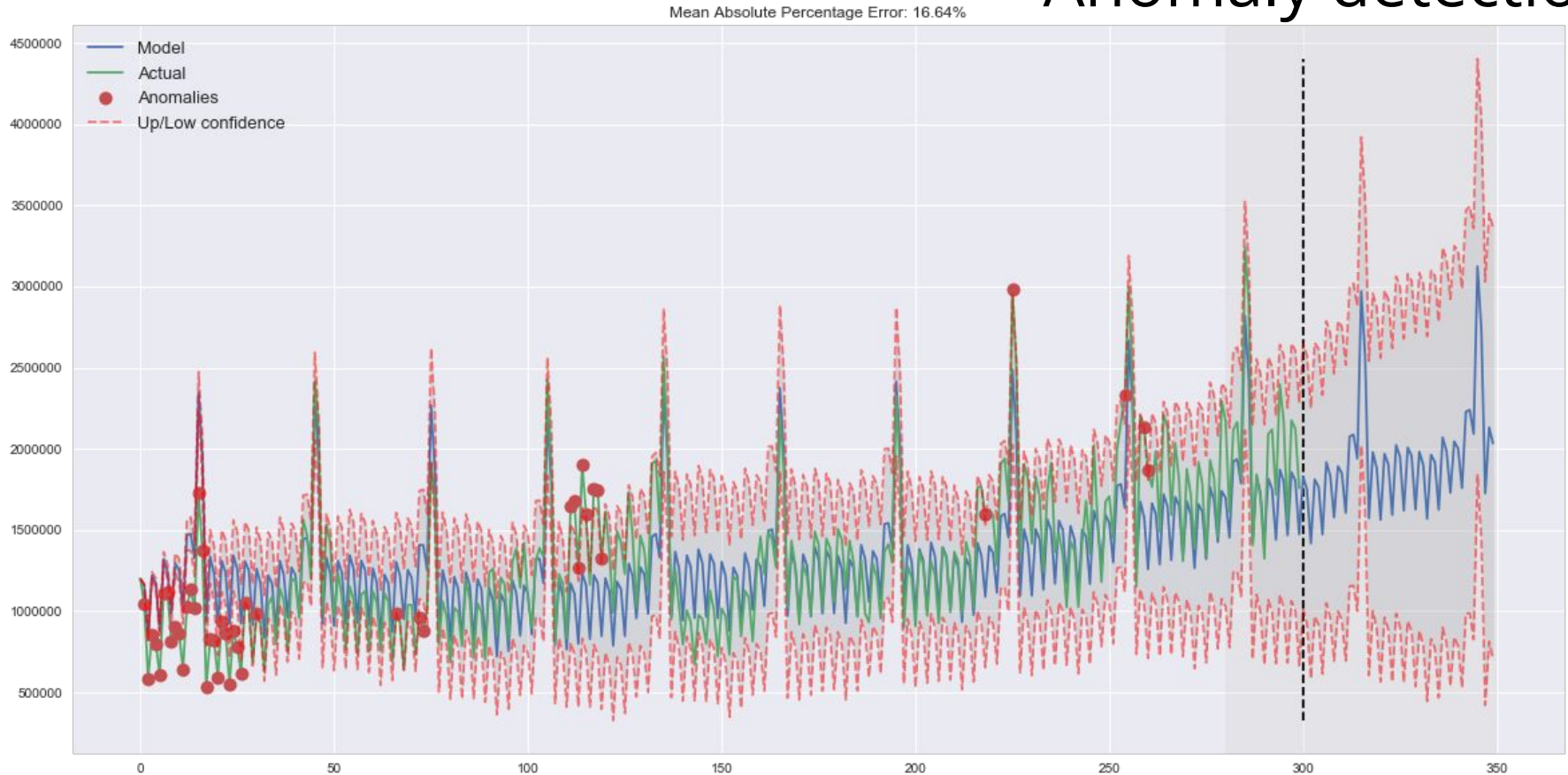
$$\hat{y}_{x+m} = \ell_x + mb_x + s_{x-L+1+(m-1)modL}$$

Mean Absolute Percentage Error: 16.64%
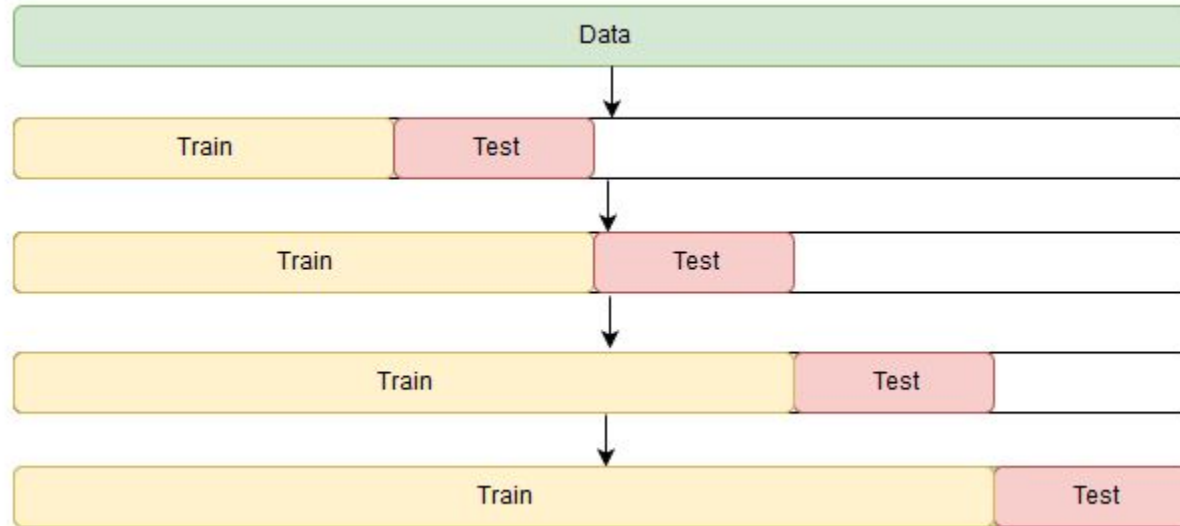
# Triple exponential smoothing example: Anomaly detection

Actually, simple exponential smoothing allows to detect anomalies as well.

One more great approach from econometrics: ARIMA (link1, link2, link3 etc.)

But first, feedback, please: http://bit.ly/ml4megafon_august18_lecture5_feedback
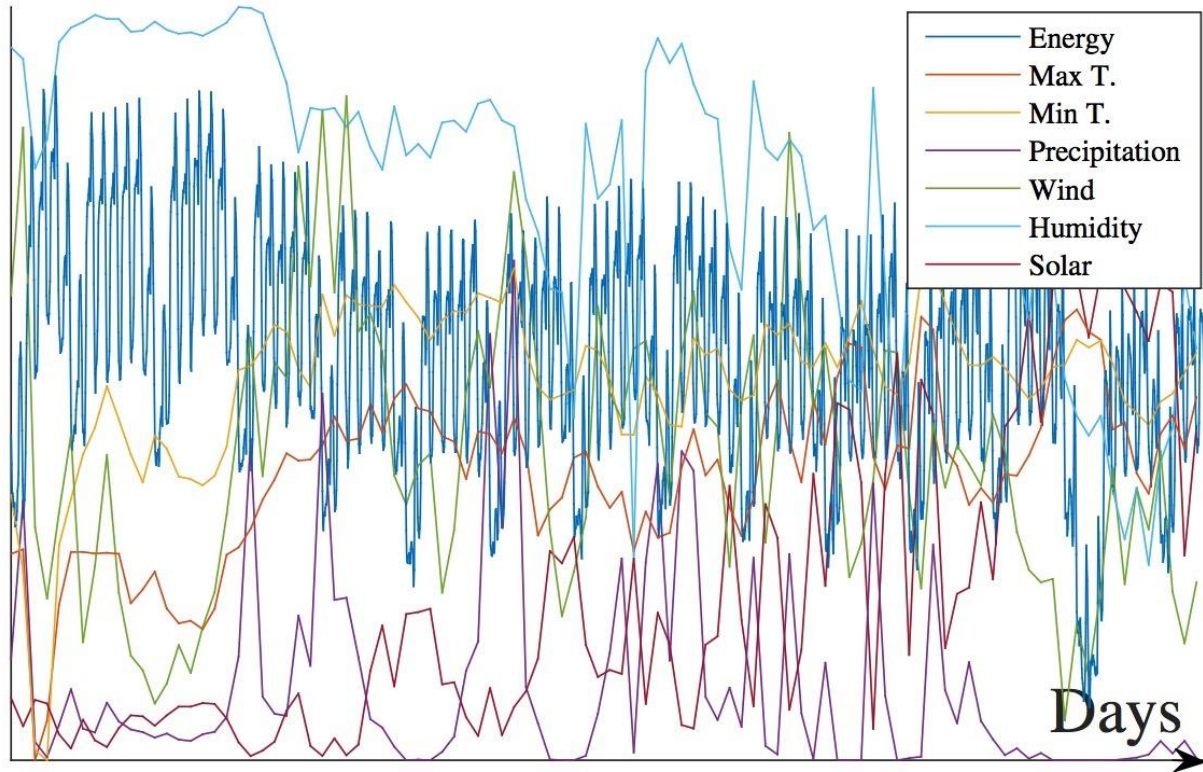
Different ways to generate features:

- Lagged Time Series Values
- Rolling statistics (max, min, median, mean, variance etc. in some window)
- Data-based features (information about holidays, siesta hours, special events etc)
- Exogenous variables
- Predictions of some extra models

# Example TS with exogenous variables
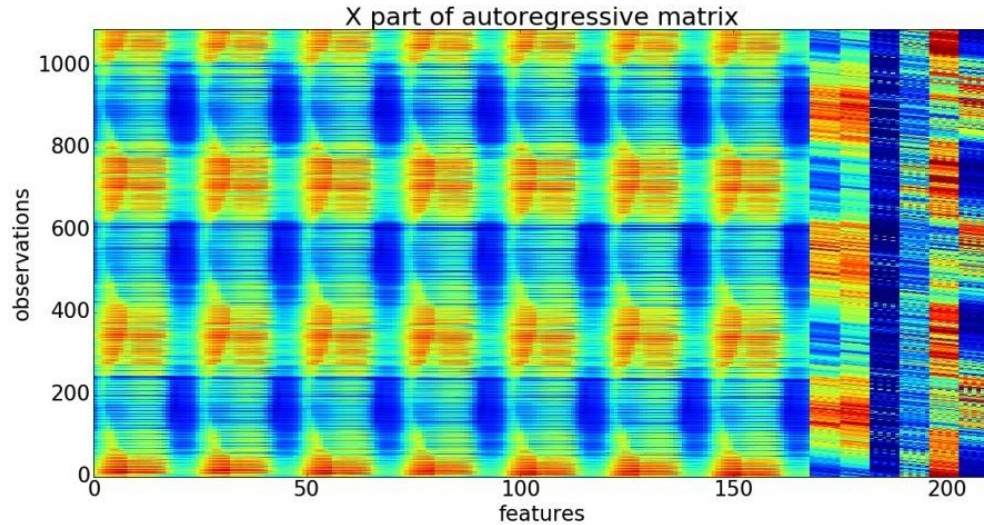
Wait, what about classification problem?
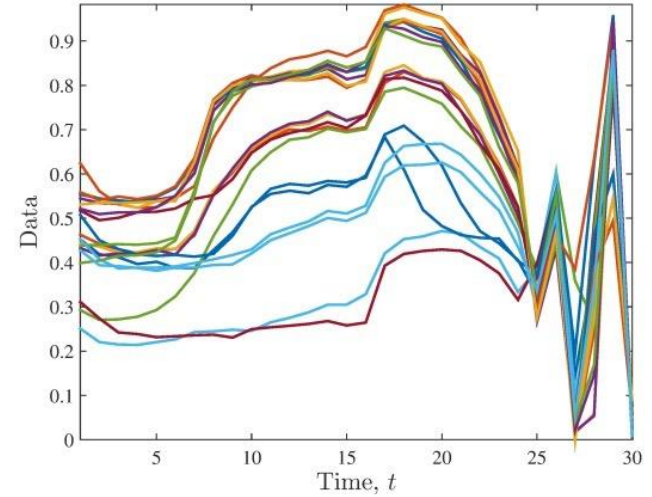
Wait, what about classification problem?

Easy: now we are just dealing with supervised learning problem!

a)



b)

a) Design matrix example
b) Target variables example

And here come all the great methods:

- Linear models
- Random Forest
- XGBoost (and alike)
- Neural Networks (RNN, CNN etc.)

To describe our data well, we need good feature space.

And feature engineering is an *art*.

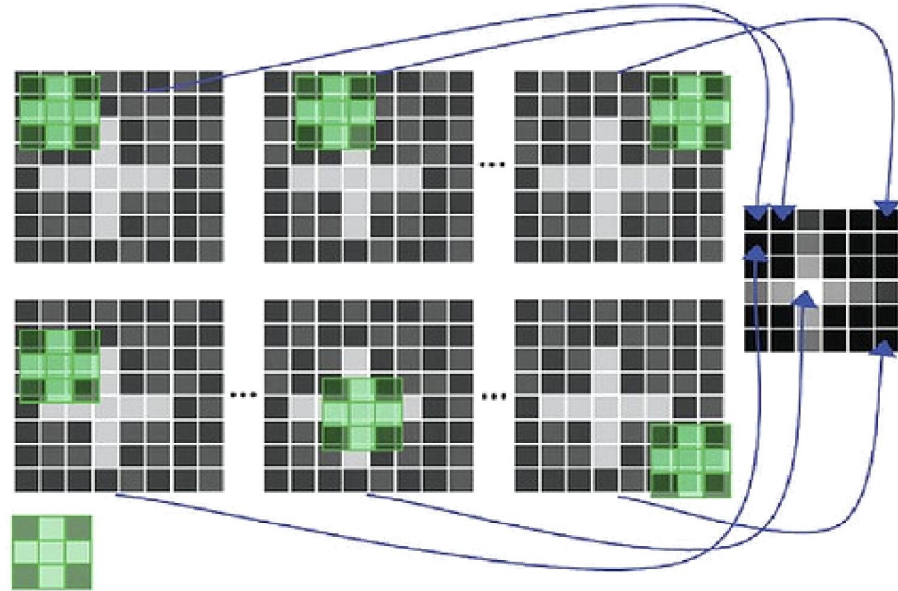Possible way out:

1. Neural Networks

2d (and 3d) convolution: vital for Computer Vision problems.

1d convolution fits time series analysis.

5x5



3x3 (5-3+1)

Convolved Feature

Image

CONVOLUTIONS
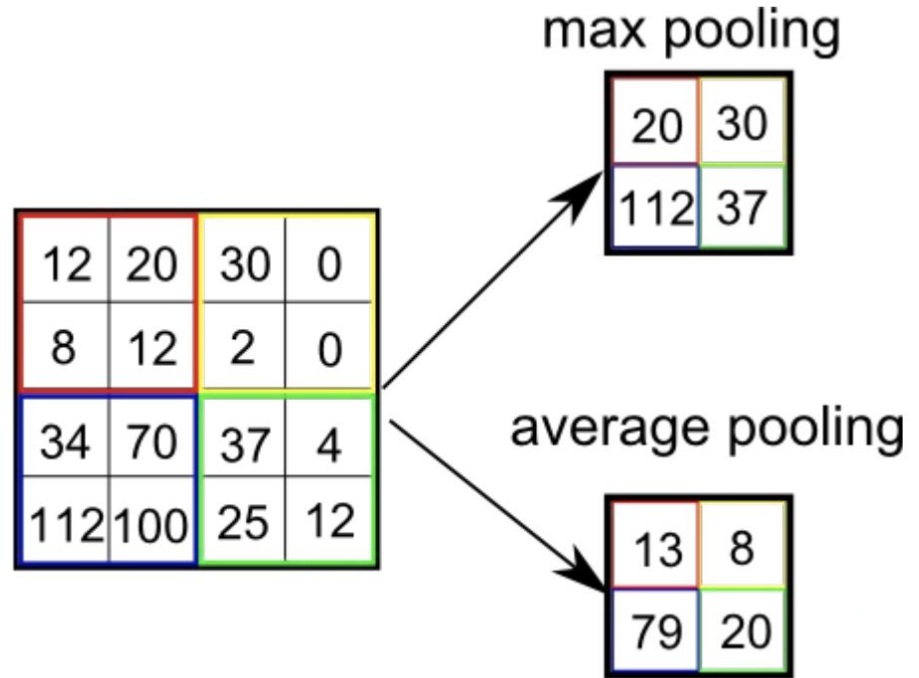
Intuition: how *cat-like* is this square?
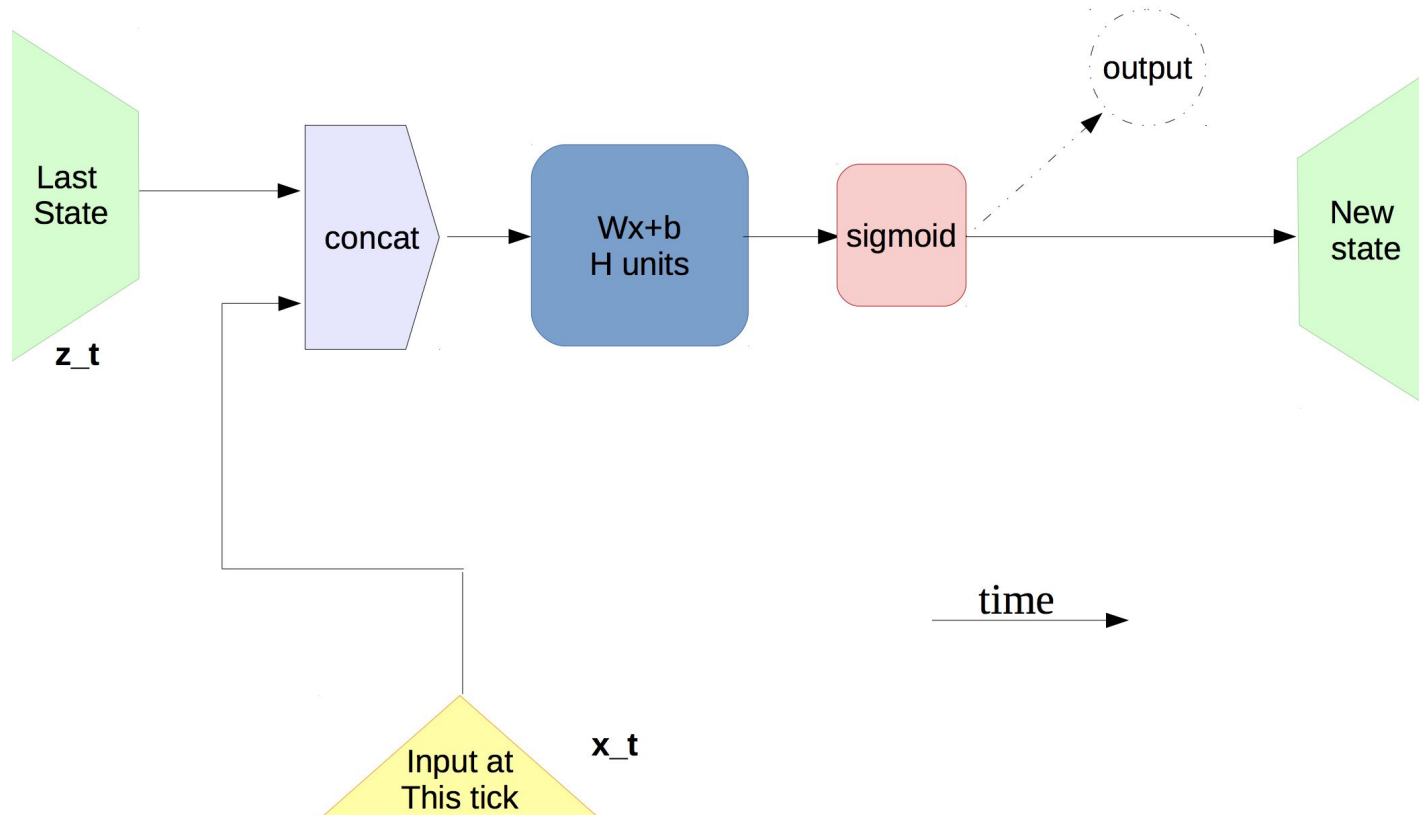
- Reduces layer size by a facto
- Makes NN less sensitive to small signal (image/ts) shifts
- 
- Widely used:
  - max pooling
  - mean pooling



max pooling

| 20 | 30 |
|----|----|
| 112 | 37 |

average pooling

| 13 | 8 |
|----|----|
| 79 | 20 |

| 12 | 20 | 30 | 0 |
|----|----|----|---|
| 8 | 12 | 2 | 0 |
| 34 | 70 | 37 | 4 |
| 112 | 100 | 25 | 12 |

We use same weight matrices for all steps

old state h0          new state h1          newer state h2          even newer state h3

input x0              input x1              input x2

x1  x2  x3  x4

embed  embed  embed  embed

x0="the"  x1="cat"  x2="sat"  x3="on"

Now with formulas

$$h_0 = \bar{0}$$

$$h_1 = \sigma\big(\langle W_{\mathrm{hid}}[h_0, x_0]\rangle + b\big)$$

$$h_2 = \sigma\big(\langle W_{\mathrm{hid}}[h_1, x_1]\rangle + b\big) = \sigma\big(\langle W_{\mathrm{hid}}[\sigma\big(\langle W_{\mathrm{hid}}[h_0, x_0]\rangle + b, x_1]\rangle + b\big)$$

$$h_{i+1} = \sigma\big(\langle W_{\mathrm{hid}}[h_i, x_i]\rangle + b\big)$$

$$P(x_{i+1}) = \mathrm{softmax}\big(\langle W_{\mathrm{out}}, h_i\rangle + b_{\mathrm{out}}\big)$$

To describe our data well, we need good feature space.

And feature engineering is an *art*.

Possible way out:

1. Neural Networks
2. Some libraries that do this stuff for us (Facebook Prophet e.g.). But it's coming next

# That's all. Time to get some practice.



But first, feedback, please:

http://bit.ly/ml4megafon_august18_lecture6_feedback